

Gradiente descendente (batch, stochastic e boosting)

Profs.: Eduardo Vargas Ferreira
Walmes Marques Zeviani

Solução de quadrados mínimos

- Seja $\mathbf{X} \in M_{n \times p}(\mathbb{R})$, com $n > p$ e $\text{posto}(\mathbf{X}) = p$. Dado $\mathbf{y} \in \mathbb{R}^n$, definimos o seguinte problema de minimização:

$$\|\mathbf{X}\hat{\boldsymbol{\beta}} - \mathbf{y}\|_2^2 = \min \left\{ \|\mathbf{X}\boldsymbol{\beta} - \mathbf{y}\|_2^2 ; \boldsymbol{\beta} \in \mathbb{R}^p \right\}$$

- Dizemos que o elemento $\hat{\boldsymbol{\beta}}$ é uma solução de quadrados mínimos;

Teorema: Seja $\mathbf{X} \in M_{n \times p}(\mathbb{R})$, com $n > p$ e $\text{posto}(\mathbf{X}) = p$. Definimos $J : \mathbb{R}^p \rightarrow \mathbb{R}$ da seguinte forma:

$$J(\boldsymbol{\beta}) = \langle \mathbf{X}\boldsymbol{\beta} - \mathbf{y}, \mathbf{X}\boldsymbol{\beta} - \mathbf{y} \rangle ; \boldsymbol{\beta} \in \mathbb{R}^p.$$

Então, o **Problema de Minimização**: encontrar $\hat{\boldsymbol{\beta}} \in \mathbb{R}^p$ tal que

$$J(\hat{\boldsymbol{\beta}}) = \min \{ J(\boldsymbol{\beta}) ; \boldsymbol{\beta} \in \mathbb{R}^p \}$$

é equivalente ao **Sistema Normal**

$$\mathbf{X}^t \mathbf{X} \boldsymbol{\beta} = \mathbf{X}^t \mathbf{y}.$$

Método do gradiente descendente (GD)

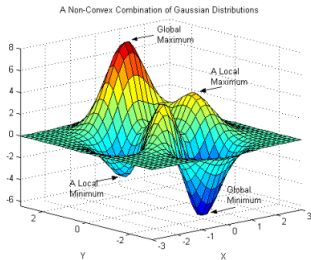


- O **Gradiente descendente** (GD) é um método para encontrar o mínimo de uma função de forma iterativa;

Algoritmo: Escolha um chute inicial, $\beta^{(0)} \in \mathbb{R}^p$, repita:

$$\beta^{(k+1)} = \beta^{(k)} - \alpha_k \nabla J(\beta^{(k)}), \quad k = 0, 1, \dots$$

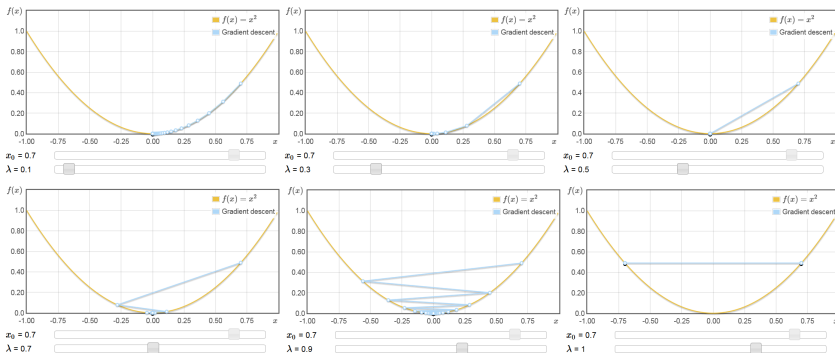
pare quando atingir convergência.



Taxa de aprendizagem α



- Taxa de aprendizagem controla o tamanho do passo em cada iteração;
- Selecionar o valor correto é crítico
 - ★ Se tomarmos α pequeno, o método fica lento;
 - ★ Se α muito grande, o método diverge.



- ✓ Ideia simples e cada iteração é barata;
 - ✓ Garantia de convergência para o mínimo local;
 - ✓ Com vários algoritmos de segunda ordem para acelerar sua convergência;
 - ✓ Muito rápido para matrizes bem condicionadas e problemas fortemente convexos;
 - ✗ Frequentemente é lento, pois problemas interessantes não são fortemente convexos ou bem condicionados;
 - ✗ Não lida com funções não diferenciáveis (dica: use o método Subgradiente).
 - ✗ Utiliza todos os dados de treinamento para estimar os parâmetros. Assim, para grandes bancos de dados, torna-se lento;
- Diante deste último aspecto, por que não em cada iteração selecionar um valor na amostra, e com sua informação executar um passo?

Gradiente descendente estocástico (GDE)

- Como vimos, no Gradiente Descendente utilizamos a **amostra completa** para atualizar os parâmetros (é um processo determinístico);
- Assim, se o tamanho da amostra de treino for grande (na verdade MUITO grande!) o Gradiente Descendente levará muito tempo em cada passo;
- A diferença no **Gradiente Descendente Estocástico** (GDE) está na utilização de somente uma observação em cada iteração;
- Então, cada passo é realizado com uma v.a. de um processo estocástico. Tornando o método mais atrativo;
- Em redes neurais, p.ex., o custo para se fazer **backpropagation** com os dados completos é alto, precisando de abordagens estocásticas como esta.

- Considere o par (x_i, y_i) amostrado do treinamento. A atualização dos parâmetros é dada por

Algoritmo: Escolha um chute inicial, $\beta^{(0)} \in \mathbb{R}^{p+1}$, repita:

$$\beta^{(k+1)} = \beta^{(k)} - \alpha_k \nabla J(\beta^{(k)}; x_i, y_i), \quad k = 0, 1, \dots$$

pare quando atingir convergência.

- No GDE a taxa de aprendizagem, α , é, tipicamente, menor do que o GD (*batch*). Isso ocorre, pois temos uma maior variância nas atualizações;
- Uma escolha de α , que funciona bem na prática, é uma taxa pequena o suficiente que dê uma convergência estável nas iterações iniciais;
- Métodos mais sofisticados incluem o uso de *Backtracking line search* ou *Exact line search*.

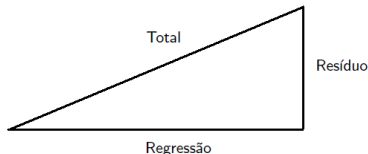
- ✓ Convergência mais rápida, especialmente com grandes bancos de dados ou dados redundantes, p. ex.:
 - Imagine que temos dados de treino de tamanho 100.000;
 - Mas na verdade são 100 cópias de 1000;
 - Ou seja, padrões parecidos, com mesmo efeito;
 - Batch será, pelo menos, 100 vezes mais devagar.
- ✓ A trajetória estocástica permite escapar de um mínimo local;
- ✗ Prova da convergência é probabilística;
- ✗ Muitos métodos de segunda ordem não funcionam;

Gradiente boosting

- Lembrando de regressão

$$\underbrace{\sum_{i=1}^n (y_i - \bar{y})^2}_{SQT} = \underbrace{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}_{SQR} + \underbrace{\sum_{i=1}^n (y_i - \hat{y}_i)^2}_{SQE}.$$

- E, geometricamente, temos

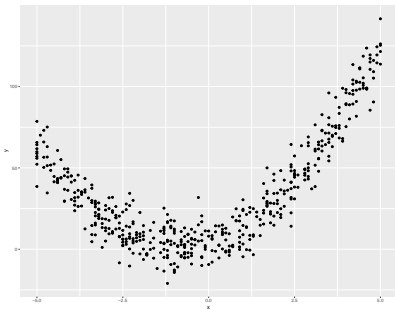


- Isso quer dizer que toda variabilidade **não explicada** pela regressão ficará no resíduo (variáveis e funções delas!);
- Vejamos um exemplo.

- Simular uma situação na qual a **verdadeira** relação entre X e Y é

$$y = 3,5x^2 + 6x + 5$$

```
x <- sample(seq(from = -5, to = 5, by = 0.1), size = 500, replace = TRUE)
y <- 3.5*x*x + 6*x + 5 + rnorm(500,0,10)
```



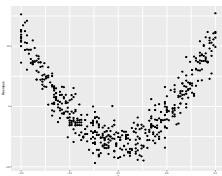
- Agora, vamos ajustar um modelo de regressão linear simples

```
ajuste = lm(y ~ x)
ajuste$coef
```

```
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  33.3499     1.2533   26.61  <2e-16 ***
## x            6.2182     0.4409   14.11  <2e-16 ***
```

- **Pergunta:** Se o termo quadrático não está no modelo, onde ele estará?
- **Resposta:** Nos resíduos.

```
e = ajuste$residuals
qplot(x,e,
      xlab="x",
      ylab="Valores ajustados")
```



- Então, considere o seguinte procedimento

$$Y = h(x) + \textit{residuo} \quad (1)$$

- Se o *residuo* não for um ruído branco (mas algo correlacionado com Y)

$$\textit{residuo} = g(x) + \textit{residuo2} \quad (2)$$

- Combinando (1) e (2)

$$Y = h(x) + g(x) + \textit{residuo2}$$

- Pode-se dizer que $h(x)$ foi atualizada com uma parte do *residuo*, ou seja

$$h(x)^{(2)} = h(x)^{(1)} + g(x)$$

- Mas, como isto está relacionado com **Gradiente boosting**?

Gradiente boosting e resíduos

- Queremos minimizar

$$J(y_i, h(\mathbf{x})) = \frac{1}{2n} \sum_{i=1}^n [y_i - h(\mathbf{x}_i)]^2$$

- Derivando com relação a $h(\mathbf{x}_i)$ temos

$$\frac{\partial J(y_i, h(\mathbf{x}))}{\partial h(\mathbf{x}_i)} = h(\mathbf{x}_i) - y_i.$$

- Podemos interpretar os resíduos como o negativo do gradiente

$$\text{resíduos} = y_i - h(\mathbf{x}_i) = -\frac{\partial J(y_i, h(\mathbf{x}))}{\partial h(\mathbf{x}_i)}$$

- Então, considerando perda quadrática, concluímos que

resíduo \Leftrightarrow negativo do gradiente
Atualizar $h(\mathbf{x}_i)$ com o resíduo \Leftrightarrow Atualizar $h(\mathbf{x}_i)$ com o negativo do gradiente

Algoritmo: Escolha um chute inicial, $h(\mathbf{x}_i)^{(0)}$, faça:

* Calcule $-\frac{\partial J(y_i, h(\mathbf{x})^{(k)})}{\partial h(\mathbf{x}_i)^{(k)}}$;

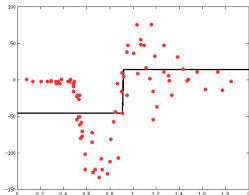
* Ajuste um modelo de regressão $g(\mathbf{x}_i)^{(k)}$ baseado no negativo do gradiente;

$$h(\mathbf{x}_i)^{(k+1)} = h(\mathbf{x}_i)^{(k)} + \rho g(\mathbf{x}_i)^{(k)}, \quad k = 0, 1, \dots$$

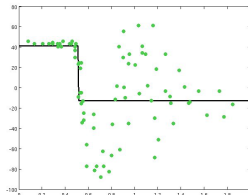
pare quando atingir convergência.

Exemplo:

Começando com um simples preditor



Aprimorando com os resíduos



Algoritmo: Escolha um chute inicial, $h(\mathbf{x}_i)^{(0)}$, faça:

* Calcule $-\frac{\partial J(y_i, h(\mathbf{x})^{(k)})}{\partial h(\mathbf{x}_i)^{(k)}}$;

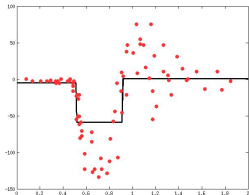
* Ajuste um modelo de regressão $g(\mathbf{x}_i)^{(k)}$ baseado no negativo do gradiente;

$$h(\mathbf{x}_i)^{(k+1)} = h(\mathbf{x}_i)^{(k)} + \rho g(\mathbf{x}_i)^{(k)}, \quad k = 0, 1, \dots$$

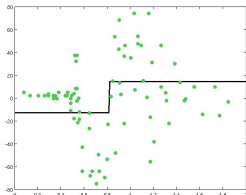
pare quando atingir convergência.

Exemplo:

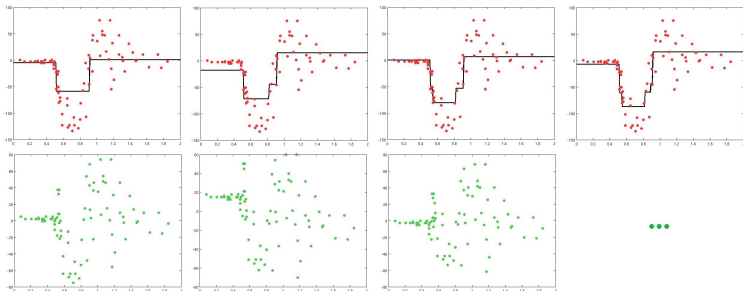
Combinando, temos um melhor preditor



Novamente, aprimorando com os resíduos



- O princípio básico é esse; propor um modelo e aprimorá-lo (ou “ensiná-lo”) através da análise dos resíduos;



- Note que podemos considerar outras funções perda e derivar o algoritmo da mesma maneira.

- **Soma dos desvios absolutos (SDA)**

$$J(y_i, h(\mathbf{x})) = \frac{1}{n} \sum_{i=1}^n |y_i - h(\mathbf{x}_i)|$$

★ O negativo do gradiente fica

$$-\frac{\partial J(y_i, h(\mathbf{x}))}{\partial h(\mathbf{x}_i)} = \text{sign}(y_i - h(\mathbf{x}_i)) = \begin{cases} 1, & \text{se } |y_i - h(\mathbf{x}_i)| < 0, \\ -1, & \text{se } |y_i - h(\mathbf{x}_i)| > 0 \end{cases}$$

- **Huber-M cost**

$$J(y_i, h(\mathbf{x})) = \frac{1}{n} \sum_{i=1}^n \begin{cases} \frac{1}{2} [y_i - h(\mathbf{x}_i)]^2, & \text{para } |y - h(\mathbf{x}_i)| \leq \delta, \\ \delta |y_i - h(\mathbf{x}_i)| - \frac{1}{2} \delta^2, & \text{caso contrário.} \end{cases}$$

★ O negativo do gradiente fica

$$-\frac{\partial J(y_i, h(\mathbf{x}))}{\partial h(\mathbf{x}_i)} = \begin{cases} y_i - h(\mathbf{x}_i), & \text{se } |y_i - h(\mathbf{x}_i)| \leq \delta, \\ \delta \text{sign}(y_i - h(\mathbf{x}_i)), & \text{caso contrário.} \end{cases}$$

- O método introduz um novo modelo de regressão em cada iteração, a fim de compensar as deficiências do modelo existente;
- As deficiências são identificadas pelo negativo do gradiente;
- Para qualquer função perda podemos derivar o Gradiente boosting;
- Perda absoluta e Huber são mais robustos a outliers;
- Para detalhes de como escolher o valor de δ , veja

▶ Greedy Function Approximation: A Gradient Boosting Machine